

Draft Report for Review Purposes Only

COMPARISON OF NATIONAL AND REGIONAL SEDIMENT QUALITY GUIDELINES FOR PREDICTING SEDIMENT TOXICITY IN CALIFORNIA

Draft Final Report

**Steven M. Bay¹, Kerry J. Ritter¹, Doris E. Vidal-Dorsch¹, and L. Jay
Field²**

**¹Southern California Coastal Water Research Project
3535 Harbor Blvd., Suite 110, Costa Mesa, CA 92626
www.sccwrp.org**

²National Oceanic and Atmospheric Administration, Seattle, WA

October 24, 2007

Technical Report

EXECUTIVE SUMMARY

A number of sediment quality guidelines (SQGs) have been developed for relating chemical concentrations in sediment to their potential for biological effects, but there have been few studies evaluating the relative effectiveness of different SQG approaches. Here we apply six SQG approaches to assess how well they predict toxicity in California sediments. Four of the SQG approaches were nationally derived indices that were established in previous studies (ERM, LRM, SQGQ1, Consensus), and two were variations of nationally derived approaches that were recalibrated to California-specific data (CA LRM and CA ERM). Each SQG approach was applied to a standardized set of matched chemistry and toxicity data for California and an index of the aggregate magnitude of contamination (e.g., mean SQG quotient or maximum probability of toxicity) was calculated. A set of three thresholds for classification of the results into four categories of predicted toxicity was established for each SQG approach using a statistical optimization procedure. The performance of each SQG approach was evaluated in terms of correlation and categorical classification accuracy. The CA LRM had the best overall performance, but the magnitude of differences in classification accuracy among the SQG approaches was relatively small. Recalibrating the indices using California data improved performance of the LRM, but not the ERM. The LRM approach is more amenable to revision than other national SQGs, which is a desirable attribute for use in programs where the ability to incorporate new information or chemicals of concern is important. As the differences in performance among indices were generally small, characteristics such as ease of application, types of chemicals included in the constituent array, and feasibility for revision, become important considerations when selecting a preferred SQG approach.

Draft Report for Review Purposes Only

ACKNOWLEDGMENTS

We thank Chris Beegan from the California Water Resources Control Board, and Mike Connor and Bruce Thompson of the San Francisco Estuary Institute for their suggestions on the design of this study. Peggy Myre of Exa Data and Mapping compiled and standardized the data sets. Jeff Brown, Diana Young, and Darrin Greenstein assisted with data compilation and statistical analysis. We also thank Peter Landrum, Ed Long, Todd Bridges, Tom Gries, Rob Burgess and Bob Van Dolah for their thoughtful review of the ideas contained within the document.

Work on this project was funded by the California State Water Resources Control Board under agreement 01-274-250-0.

INTRODUCTION

Many monitoring programs are conducted to evaluate chemical contamination effects on sediment quality, but interpreting these data is difficult (Wenning *et al.* 2005). The biological availability of chemicals in sediments is complex and poorly understood. Moreover, the chemicals are often present in complex mixtures that are difficult to integrate.

A number of sediment quality guidelines (SQGs) for relating chemical concentrations to potential for biological effects have been developed, generally falling into two classes. The first is a mechanistic approach, which models the chemical and biological processes that affect contaminant bioavailability. Current mechanistic SQGs are based on equilibrium partitioning theory and apply to selected classes of contaminants, primarily divalent metals and some types of nonionic organics (USEPA 2004a, 2004b). While these models are useful for describing potential contaminant availability and identifying the cause of toxicity, mechanistic SQGs are not available for many contaminants of interest and they correlate poorly with biological effects under field conditions (Vidal and Bay 2005). In addition, some of the parameters needed to apply these guidelines (e.g., sediment acid volatile sulfides and simultaneously extracted metals) are rarely collected in routine monitoring programs.

A more widely used type is empirical SQGs, which are guidelines derived from the statistical analysis of matched sediment chemistry and biological effects data. Multiple collections of empirical SQGs that are based on different statistical approaches have been developed. Examples of empirical SQG approaches for the marine environment include the effects range-median (ERM), probable effects level (PEL), apparent effects threshold (AET), sediment quality guideline quotient (SQGQ1), and logistic regression models (LRM) (Barrick *et al.* 1988, Fairey *et al.* 2001, Field *et al.* 2002, Long *et al.* 1995, MacDonald *et al.* 1996). Consensus guidelines, which aggregate several different of SQGs having a similar narrative intent (e.g., median effect), are an evolution of the empirical approach. Marine consensus SQGs have been developed for a some constituents, including metals, PCBs, and PAHs (MacDonald *et al.* 2000, Swartz 1999, Vidal and Bay 2005).

It is unclear which empirical SQG approach is most effective for describing the potential for biological effects associated with chemical contamination. Numerous studies have shown that each SQG approach has predictive ability with respect to biological effects, but most studies have generally been limited to examination of just one or two approaches and often use variable methods to measure performance (Wenning *et al.* 2005). Long *et al.* (2000) applied ERMs and PELs to several data sets and observed different patterns in predictive ability. Vidal and Bay (2005) compared five SQG approaches using a common data set and found large differences in predictive ability among some approaches, however their study did not include the logistic regression approach. Vidal and Bay (2005) also observed that comparisons of SQG performance can be strongly influenced by the selection of thresholds used to classify the results. Existing studies are inadequate for comparing the performance of empirical SQGs because of their limited scope, lack of comparability in methods, and lack of thresholds derived using a consistent methodology.

Draft Report for Review Purposes Only

It is also unclear whether performance of SQGs is improved when they are calibrated to local conditions. The predictive ability of SQGs to biological effects has been shown to vary when the same guidelines are applied to data from different regions (Fairey *et al.* 2001, Long *et al.* 1998, Long *et al.* 2006, O'Connor *et al.* 1998, Vidal and Bay 2005). These variations in performance may be due to differences in the nature of the chemical mixtures between sites or regions, variations in bioavailability due to geochemical factors, or differences in the sensitivity of methods used to measure biological effects. Variation in SQG performance among studies creates uncertainty in determining the threshold of SQG exceedance associated with adverse impacts on sediment quality. The use of SQGs and interpretation thresholds that are derived or calibrated relative to site-specific conditions has been recommended as a way to reduce the uncertainty of SQG interpretation (Fairey *et al.* 2001, Long *et al.* 2006, Vidal and Bay 2005).

Here we apply six SQG approaches to a large California data set of paired chemistry and toxicity measurements to assess: 1) which SQG approach best predicts toxicity of California sediments, 2) whether the ability of SQGs to predict sediment toxicity is improved when the SQGs are recalibrated to California data, and 3) if performance further improves when the SQGs are further recalibrated to two subregions within California.

METHODS

We assessed the performance of six SQG approaches by applying them to matched chemistry and toxicity data for California, calculating an index of overall contamination based on the mean SQG quotient or the maximum probability of toxicity, and determining the correlation and categorical classification accuracy. Four of the SQG approaches were derived in previous national studies (ERM, LRM, SQGQ1, Consensus) and two were variations of nationally derived SQGs that were recalibrated to California-specific data (CA LRM and CA ERM). SQG calibration and performance evaluations were conducted at two scales in order to investigate the influence of regional variations in sediment characteristics: statewide (all California data) and regional (separate northern and southern California data sets).

Data

Paired chemistry and sediment toxicity measurements from California marine embayments were compiled from 151 dredging, monitoring, and research studies conducted in California between 1984 and 2004. The database included stations from marine and estuarine embayments located from 41.94°N (Del Norte County, CA, USA) to 31.75°N (USA-Mexico international border). More information on the studies used to populate this database can be found at http://www.sccwrp.org/data/2006_sqg.html.

The data were screened to select information that was of high quality and comparable. All stations were from locations in enclosed bays or harbors at subtidal depths and only data from surficial sediment (top 30 cm or less) were selected. Toxicity data were limited to information from solid-phase 10-d amphipod survival tests using *Rhepoxynius abronius* or *Eohaustorius estuarius* and conducted using standardized methods (USEPA 1994). The toxicity data were further screened to ensure that conventional data quality objectives were met, including mean control survival >85% and overlying water ammonia concentrations below species-specific criteria (USEPA 1994). Screening steps to select chemistry data for analysis included a review of the data quality assessment from the study authors, use of comparable extraction/digestion methods, and measurement of a minimum suite of contaminants that included multiple metals and PAHs.

Standardized sums of PAHs, DDTs, PCBs, and chlordanes were calculated using a consistent methodology for all samples. Low molecular weight PAHs were calculated as the sum of acenaphthene, anthracene, biphenyl, naphthalene, 2,6-dimethylnaphthalene, fluorene, 1-methylnaphthalene, 2-methylnaphthalene, 1-methylphenanthrene, and phenanthrene. High molecular weight PAHs was the sum of benzo(a)anthracene, benzo(a)pyrene, benzo(e)pyrene, chrysene, dibenz(a,h)anthracene, fluoranthene, perylene, and pyrene. Total PAHs was the sum of Low PAH and the High PAH values. Total PCBs were calculated from the sum of congeners 8, 18, 028, 44, 52, 66, 101, 105, 110, 118, 128, 138, 153, 180, 187, and 195. This sum was multiplied by 1.72 to estimate the total concentration of all congeners. Total DDTs represented the sum of p,p'-DDT, o,p'-DDT, p,p'-DDE, o,p'-DDE, p,p'-DDD, and o,p'-DDD. Total chlordanes was the sum of alpha-chlordane (cis-chlordane), oxychlordane, trans-chlordane, trans-nonachlor, and gamma-chlordane.

Draft Report for Review Purposes Only

Data were estimated for values reported as below reporting limits based on multiple regression imputation, taking advantage of covariation among the many chemical and sediment variables. Imputation produces lesser bias than conventional approaches for interpreting nondetect data, such as substituting zero or one-half of the reporting limit (Helsel 2005). SAS PROC MI (SAS Institute Inc, North Carolina, USA) was used to impute values in a sequential stepwise fashion by contaminant type. Metal data were estimated first, followed in order by pesticides, PAHs, and PCBs. The stepwise manner in which the groups of data variables were imputed was used because SAS PROC MI could not compute all imputations in a single step. The stepwise procedure also allowed for better control of the data variables used in the imputations for each chemical group.

The standardized data set was divided into two groups to facilitate investigation of regional differences in chemical contamination on SQG performance: northern California embayments north of Pt. Conception and southern California embayments south of Pt. Conception. Each regional data set was further divided into two portions: a calibration subset used for index development and threshold calibration, and an independent validation subset used for the analysis of SQG performance. Approximately one third of the data were used for validation. The validation samples were selected by first grouping the data into one of 8 subregions based on latitude to ensure even spatial representation. The samples within each subregion were then ranked by the mean mERMq quotient and one third of the samples systematically sampled from throughout the mERMq quotient distribution. Additional validation data were obtained from recent monitoring studies that were not included in the initial data compilation effort. The north and south validation data sets contained 146 and 249 samples, respectively.

National SQGs

The ERM values used in the analyses were obtained from Long *et al.* (1995). The mean ERM quotient (mERMq) for each sample in the data set was calculated by dividing each chemical concentration by its respective ERM and averaging the individual quotients (Long *et al.* 2000). The subset of ERM values used to calculate the mERMq (Table 1) was the same as that used in previous mERMq performance studies (Long *et al.* 2000).

The mean sediment quality guideline quotient 1 (SQGQ1) was calculated as described by Fairey *et al.* (2001). The SQG values used in the analysis are listed in Table 1.

The Consensus SQG values for PAHs and PCBs were midpoint effect concentrations obtained from Swartz (1999) and MacDonald *et al.* (2000), respectively. Values for DDTs, dieldrin, arsenic, cadmium, chromium, copper, lead, mercury, nickel, silver, and zinc were obtained from Vidal and Bay (2005). The mean Consensus quotient was calculated by dividing each chemical concentration by its respective SQG (Table 1) and averaging the individual quotients.

The Logistic Regression Model (LRM) approach was based on the statistical analysis of paired chemistry and amphipod toxicity data from studies throughout the U.S. (Field *et al.* 1999, 2002). The logistic model is described by the following equation:

$$p = e^{B_0 + B_1(x)} / (1 + e^{B_0 + B_1(x)})$$

Draft Report for Review Purposes Only

where: p = probability of observing a toxic effect;
B0= intercept parameter;
B1= slope parameter; and,
x= concentration or log concentration of the chemical.

The chemical-specific models used in this study were based on an analysis of the accuracy for predicting toxicity for 37 candidate models. Models for 18 chemicals having low rates of false positives were selected for use (Table 2). The maximum probability of toxicity obtained from the individual models (P_{\max}) for each sample was used as the index of overall contamination.

Regional SQGs

Regionally calibrated versions were developed for two of the national SQG approaches: ERM and LRM. Regional versions were not developed for the other national SQG approaches (SQGQ1 and Consensus) because these approaches are based on the inclusion of SQG values from other sources and cannot be easily recalibrated with new data. Three versions of each regional SQG approach were developed: a statewide version that was calibrated to data from throughout California (CA ERM or CA LRM), and two region-specific versions. The region-specific versions were calibrated separately for the northern California (NorCA ERM or NorCA LRM) and southern California (SoCA ERM or SoCA LRM) data sets.

For the CA ERM variations, local calibration involved calculation of new individual chemical ERM values. The data were screened to select toxic samples (>20% mortality) with chemical concentrations >2x median concentration of nontoxic samples. A separate screening process was used for each chemical. After screening, the data were sorted in ascending order and the median concentration of each chemical was selected as the region-specific ERM value. ERM values were calculated for all chemicals having >10 records in the screened data set. This resulted in calculating CA ERM and SoCA ERM values for 27 chemicals, and NorCA ERM values for 25 chemicals (Table 1).

California logistic regression models for individual chemicals were developed for the statewide and regional California data sets using the methods described in USEPA (2005). These models were applied to the California calibration data using <80% control-adjusted amphipod survival as the definition of a toxic sample. The specific models included in the CA LRM, SoCA LRM, and NorCA LRM approaches were selected from a library of candidate models that included national models, as well as models derived using the California data sets. The selected models were chosen based on the goodness of fit with the observed probability of toxicity (Table 2). Models with high false positive rates were not included.

Threshold Development

Evaluating the indices with respect to categorical classification accuracy requires identification of category thresholds for each SQG index. Such thresholds are generally unavailable for these SQG approaches or vary in the method of development. The thresholds used in this study were developed for each SQG approach using a consistent methodology so that differences in performance would reflect inherent differences among approaches, rather than variations in how thresholds were assigned.

Draft Report for Review Purposes Only

Three thresholds, defining four ranges of SQG index results, were established for each SQG approach. Each SQG index range corresponded to one of four categories of toxicological response that were based on classification systems used in other studies (Long *et al.* 2006). The toxicity categories were specific to each test species and were based on analyses of the minimum significant difference and magnitude of response (percent of control survival) to California samples (Bay *et al.* 2007). The categories for *E. estuarius* were: nontoxic ($\geq 90\%$ survival), low toxicity (82-89%), moderate toxicity (59-81%) and high toxicity ($< 59\%$). The categories for *R. abronius* were: nontoxic ($> 90\%$ survival), low toxicity (83-89%), moderate toxicity (70-82%) and high toxicity ($< 70\%$).

The thresholds were selected using a statistical optimization procedure based on maximizing overall agreement between the SQG index and toxicity categories in the calibration data set. The percent agreement was computed for all possible sets of triplicate thresholds occurring within a relatively dense set of possibilities. Mesh sizes for optimization reflected a distance between possible thresholds values of 5% of the range of data values for each indicator. In addition, distances between individual thresholds within each set were constrained to be no less than 10% of the range of data values for each SQG index. These constraints ensured that optimization converged and the resulting thresholds were not too close to one another. The set of triplicate thresholds that yielded the highest percent agreement were selected as being optimal.

The optimization procedure was conducted on a subset of the data that contained an even distribution of samples across toxicity categories. This step was included in order to minimize the influence on the optimization results of the skewed sample distribution in the calibration data set, which contained a higher proportion of nontoxic and low toxicity samples. The measurement of percent agreement is sensitive to skewed distributions, potentially resulting in inaccurate thresholds. The threshold selection data set contained 30 randomly selected calibration samples from each toxicity category. Data selection and threshold optimization was bootstrapped 50 times using SAS PROC SURVEYSELECT (SAS Institute Inc, North Carolina, USA) in order to provide a robust collection of thresholds that reflected variations in the calibration data. The optimum set of thresholds was determined for each iteration and the median set of thresholds was chosen to be the final thresholds for that SQG approach.

Evaluation of SQG Performance

SQG performance was evaluated by quantifying the strength of association between chemistry and toxicity in terms of both correlation and categorical classification accuracy. Correlation was measured as the nonparametric Spearman's correlation coefficient between the SQG index value (i.e., mean quotient or P_{\max}) and percent amphipod mortality. Analyses of categorical classification accuracy were based on the frequency with which the SQG index category (determined by applying the thresholds derived from the calibration data set) correctly predicted the measured toxicity response category. All analyses were conducted using an independent validation data set that was not used for threshold development.

Draft Report for Review Purposes Only

Two measures of classification accuracy were calculated: percent agreement and weighted kappa. Percent agreement is the number of samples that are correctly classified and was calculated as:

$$A=(N_c/N_t)*100$$

where: A = percent agreement
 N_c = number of samples correctly classified
 N_t = total number of samples

The weighted kappa statistic (Cohen 1960, 1968) is also measure of agreement between the SQG predictions and toxicity, but differs in that a correction for chance is applied and partial credit is given according to the severity of disagreement. Kappa weights were based on the linear weighting scheme of Cicchetti-Allsion (1971); a weight of 1 was assigned to cases of perfect agreement and weights of 1/3, 1/6, and 0 assigned to disagreements of one, two, or three toxicity categories, respectively. SAS PROC FREQ (SAS Institute Inc, North Carolina, USA) was used to calculate the weighted kappa (Stokes *et al.* 2000).

A bootstrap resampling approach similar to that used for threshold development was also used in calculation of the correlation, percent agreement, and weighted kappa values. The reported correlation and classification accuracy values are the median of 50 resamples. The 90th percentile confidence limits of the bootstrapped results were used to identify the best performing SQG approaches with respect to correlation and classification accuracy. The approach having the highest values for both correlation and classification accuracy was selected as the best performing SQG. The correlation results were given greater weight when the rankings were variable among the performance measures in order to minimize the influence of threshold selection.

Draft Report for Review Purposes Only

RESULTS

Different patterns of sediment contamination between the northern and southern California data sets (Table 3), reflecting different anthropogenic inputs and geology, were apparent. Median concentrations of most PAH compounds, chromium, and nickel were greatest in the north, while the south data set contained higher concentrations of chlordane, copper, DDTs, PCBs, and zinc. The southern California data set usually contained the highest concentrations of each contaminant, which may reflect the larger number of data in the south data set. An exception was the presence of the higher chromium and nickel concentrations in the north data set, which was likely due to higher naturally occurring concentrations of these elements in northern California soils.

There was a similar range and distribution of sediment toxicity between the northern and southern California data sets (Figure 1). The distribution of the data was skewed towards low toxicity; approximately 60% of the samples in each region had less than 20% mortality and less than 10% had greater than 60% mortality.

There were large differences in the number of chemicals and their threshold concentrations included in the different SQG indices (Tables 1 and 2). The number of chemicals varied from 9 for the SQGQ1 to 25 for the mERMQ. Individual chemical concentrations for the ERM, SQGQ1, and Consensus SQGs were similar because these values were often derived from similar sources. There were often large differences in individual chemical concentration between the national and region-specific versions of the ERM. This was especially evident for PAH compounds, where the national ERM values were 1-2 orders of magnitude greater than the CA ERMs (Table 1).

The categorization thresholds for the SQGs varied geographically (e.g., statewide, north, south). The largest thresholds were usually obtained for southern California data, but the differences were typically small (Table 4). The SQGQ1 was an exception, having nearly a three-fold difference between thresholds derived using northern and southern California data.

Each of the statewide-calibrated SQG approaches correlated significantly with amphipod survival when applied to statewide validation data. Spearman correlation coefficients ranged from 0.35 to 0.16 (Table 5), with the CA LRM having the highest correlation. Correlations generally increased when the indices were evaluated using the separate north and south data sets, though CA LRM performed best in both habitats (Table 6).

The CA LRM (Table 5) also performed best with respect to classification accuracy, when the indices were applied to the statewide data set. Very little improvement in classification accuracy was obtained using the CA ERM approach, relative to the national ERM approach. While both measures of classification accuracy ranked the SQG approaches similarly, the weighted kappa statistic provided a greater degree of discrimination among approaches than did percent agreement.

When the SQG indices and statewide thresholds were evaluated relative to the regional data sets, the CA LRM was the only approach with consistently high classification accuracy and

Draft Report for Review Purposes Only

correlations (Table 6). The CA ERM also had relatively high classification accuracy for northern California data and high classification accuracy was also obtained for the ERM and Consensus for southern California.

Developing thresholds on a regional basis had little effect overall. Percent agreement scores across indices were almost identical between thresholds developed using statewide and regional data sets (Table 6). However, classification accuracy (weighted kappa) was improved for the worst performing SQG approaches, such as SQGQ1 in the south and national LRM in the north.

Increased classification accuracy was obtained for the region-calibrated SQGs in the north (NorCA LRM and NorCA ERM) compared to statewide-calibrated versions (Table 6). However, no improvement was measured for the approaches that were calibrated to southern California data (SoCA LRM and SoCA ERM).

DISCUSSION

While the Pmax, based on the CA LRM, was the best-performing SQG index, there was relatively little difference in performance among many of the indices. This differs from the findings of Vidal and Bay (2005) and probably results from using thresholds that were selected using a consistent methodology and calibration data set. The standardized thresholds allowed each SQG approach to be evaluated on a level playing field, so that differences in performance could be compared without the confounding effect of differences in threshold selection.

Two of the SQG approaches were recalibrated using California data, which had mixed effects. For the CA LRM, there was a substantive improvement in performance, but performance of the mean quotients based on the CA ERM, was comparable to that of the national mERMQ. This may have resulted from differences in the SQG calibration process. The CA ERMs consisted of entirely of new values that were derived from the California data set. All available CA ERMs were used in the quotient calculations. In contrast, for the CA LRM, the set of models used for evaluation was selected from a combination of national and California derived models. This selection process was based on increasing model goodness of fit and reducing false positives. It is possible that this additional selection step improved the predictive ability of the CA LRM. A similar selection process was not used for the CA ERM because of differences in derivation methodology compared to the national ERMs, which were based on multiple types toxicity tests and other biological response values (Long *et al.* 1995).

The improved performance of the CA LRM may also have been due to differences in the composition, magnitude, and bioavailability of sediment contamination in the California data, relative to the data used for national LRM development. Regional differences in contamination and geochemistry have been identified as important factors affecting the predictive accuracy of SQGs (Long *et al.* 2000, Wenning *et al.* 2005). Since the values used in empirical SQG approaches are derived from chemistry-toxicity relationships in the development data set, regionally calibrated approaches would be expected to have greater predictive accuracy.

Use of thresholds calibrated to the north and south subregions produced only small increases in performance relative to the statewide thresholds. The relatively small differences in regional performance are probably related to the heterogeneous nature of sediment contamination. Even though there are differences in overall pattern and magnitude of contamination in the northern and southern California data sets, contamination patterns within each region is highly diverse due to the presence of multiple waterbodies and contaminant inputs from a multitude of sources.

Because the performance difference among SQG indices was small, characteristics such as history of use, ease of application, types of chemicals included in the constituent array, and feasibility for revision should be considered when selecting the SQG approach to be used. For instance, the Consensus and SQGQ1 approaches incorporate a lesser number of chemicals than the other approaches and it is difficult to add new contaminants of concern because the SQGs are dependent on the availability of values from other sources. Local calibration is also not feasible for these approaches for the same reason.

Draft Report for Review Purposes Only

The best performing index, CA LRM, is highly amenable to revision as demonstrated by this study. But LRM approaches are also the most difficult to apply and interpret because a complex set of regressions must be used to determine probabilities of toxicity, rather than comparing chemistry data to a simple table of SQG values. These difficulties can be overcome by incorporating the regression calculations into spreadsheets or other data analysis tools and establishing thresholds for interpreting the P_{\max} values.

Draft Report for Review Purposes Only

LITERATURE CITED

- Barrick R, Becker S, Brown L, Beller H, Pastorok R. 1988 Sediment quality values refinement: 1988 Update and evaluation of Puget Sound AET, Volume 1. PTI Environmental Services, Bellvue, WA.
- Bay S, Greenstein D, Young D. 2007. Evaluation of methods for measuring sediment toxicity in California bays and estuaries. Costa Mesa, CA: Southern California Coastal Water Research Project. Technical Report 503.
- Cicchetti, D.V. and Allison, T. 1971. A new procedure for assessing reliability of scoring EEG sleep recordings. *American Journal of EEG Technology* 11, 101-109.
- Cohen, J. 1960. A coefficient of agreement for nominal scales. *Educational and Psychological Measurement* 20:37-46.
- Cohen, J. 1968. Weighted Kappa nominal scale agreement with provision for scale disagreement or partial credit. *Psychological Bulletin* 70:213-220.
- Fairey R, Long ER, Roberts CA, Anderson BS, Phillips BM, Hunt JW, Puckett HR, CJ Wilson. 2001. An evaluation of methods for calculating mean sediment quality guideline quotients as indicators of contamination and acute toxicity to amphipods by chemical mixtures. *Environ Toxicol Chem* 20:2276-2286.
- Field LJ, MacDonald D, Norton SB, Severn CG, Ingersoll CG. 1999. Evaluating sediment chemistry and toxicity data using logistic regression modeling. *Environ Toxicol Chem* 18:1311-1322.
- Field LJ, MacDonald DD, Norton SB, Ingersoll CG, Severn CG, Smorong D, Lindskoog R. 2002. Predicting amphipod toxicity from sediments using Logistic Regression Models. *Environ Toxicol Chem* 9:1993-2005.
- Helsel D. 2005. More than obvious: better methods for interpreting nondetect data. *Environ. Sci Technol* 39:419A-423A.
- Long ER, MacDonald DD, Smith SL, Calder FD. 1995. Incidence of adverse biological effects within ranges of chemical concentrations in marine and estuarine sediments. *Environ Manag* 19:81-97.
- Long ER, Field JE, MacDonald DD. 1998. Predicting toxicity in marine sediments with numerical sediment quality guidelines. *Environ Toxicol Chem* 17:714-727.
- Long ER, MacDonald DD, Severn CG, Hong CB. 2000. Classifying the probabilities of acute toxicity in marine sediments with empirically-derived sediment quality guidelines. *Environ Toxicol Chem* 19:2598-2601.

Draft Report for Review Purposes Only

Long ER, Ingersoll CG, MacDonald DD. 2006. Calculation and uses of mean sediment quality guideline quotients: A critical review. *Environ Sci Technol* 40:1726-1736.

MacDonald DD, Carr RS, Calder FD, Long ER, Ingersoll CG. 1996. Development and evaluation of sediment quality guidelines for Florida coastal waters. *Ecotoxicology* 5:253-278.

MacDonald DD, Di Pinto LM, Field LJ, Ingersoll CG, Long ER, Swartz RC. 2000. Development and evaluation of consensus-based sediment effect concentrations for polychlorinated biphenyls (PCB). *Environ Toxicol Chem* 19:1403-1413.

O'Connor TP, Daskalakis KD, Hyland JL, Paul JF, Summers JK. 1998. Comparisons of sediment toxicity with predictions based on chemical guidelines. *Environ Toxicol Chem* 17:468-471.

Stokes, Maura E, Charles S. Davis, Gary G. Koch. 2000. *Categorical Data Analysis using the SAS system*. 2d ed. Cary NC: SAS Institute Inc.

Swartz RC. 1999. Consensus sediment quality guidelines for PAH mixtures. *Environ Toxicol Chem* 18:780-787.

[USEPA] U.S. Environmental Protection Agency. 1994. Methods for assessing the toxicity of sediment-associated contaminants with estuarine and marine amphipods. Washington DC: Office of Research and Development. EPA 600-R94-025.

[USEPA] U.S. Environmental Protection Agency. 2004a. Procedures for the derivation of equilibrium partitioning sediment benchmarks (ESBs) for the protection of benthic organisms: Metal mixtures (cadmium, copper, lead, nickel, silver, and zinc). Washington DC: Office of Research and Development. EPA-600-R-02-011.

[USEPA] U.S. Environmental Protection Agency. 2004b. Procedures for the derivation of equilibrium partitioning sediment benchmarks (ESBs) for the protection of benthic organisms: Nonionics compendium. Washington DC: Office of Research and Development. EPA-600-R-02-016.

[USEPA] U.S. Environmental Protection Agency. 2005. Predicting Toxicity to Amphipods from Sediment Chemistry (Final Report). Washington, DC, ORD National Center for Environmental Assessment. EPA/600/R-04/030.

Vidal DE, Bay SM. 2005. Comparative sediment guideline performance for predicting sediment toxicity in southern California, USA. *Environ Toxicol Chem* 24:3173-3182.

Wenning RJ, Batley GE, Ingersoll CG, Moore DW, eds. 2005. Use of sediment quality guidelines (SQGs) and related tools for the assessment of contaminated sediments. Pensacola (FL): Society of Environmental Toxicology and Chemistry.

Draft Report for Review Purposes Only

Table 1. Chemical values for individual sediment quality guidelines used for data analyses. Values for the effects range median (ERM) were taken from Long *et al.* 1995. Mean sediment quality guideline quotient (SQGQ1) values taken from Fairey *et al.* 2001. Consensus midpoint effect concentration values taken from Swartz, 1999; MacDonald *et al.* 2000; and Vidal and Bay 2005. Concentrations are on a dry weight basis except where noted.

Chemical	Units	ERM	CA ERM	SoCA ERM	NorCA ERM	SQGQ1	Consensus
Arsenic	mg/kg	70.0	19.2	19.1			55.0
Cadmium	mg/kg	9.6	1.0	1.2	0.6	4.2	5.9
Chromium	mg/kg	370.0	154.0	110	291.0		224.9
Copper	mg/kg	270.0	151.0	208	91.2	270	225.0
Lead	mg/kg	218.0	87.4	94.5	56.4	112.2	222.3
Mercury	mg/kg	0.71	0.8	0.8	0.7		0.6
Nickel	mg/kg	51.6	83.5	42			67.6
Silver	mg/kg	3.7	0.9	1.1	0.4	1.8	3.4
Zinc	mg/kg	410.0	332.5	406.9	214.5	410.0	357.1
2-Methylnaphthalene	µg/kg	670.0	22.2	23.6	20.2		
Acenaphthene	µg/kg	500.0	23.0	24.5	19.0		
Acenaphthylene	µg/kg	640.0	26.0	47	19.8		
Anthracene	µg/kg	1,100.0	130.0	215.5	60.8		
Benzo(a)anthracene	µg/kg	1,600.0	356.6	540	169.5		
Benzo(a)pyrene	µg/kg	1,600.0	405.5	630	225.3		
Chrysene	µg/kg	2,800.0	577.0	739.9	239.0		
Dibenz(a,h)anthracene	µg/kg	260.0	94.4	130	23.4		
Dieldrin	µg/kg	8.0	2.0	2	0.8	8.0	7.0
Fluoranthene	µg/kg	5,100.0	432.3	723	410.9		
Fluorene	µg/kg	540.0	30.7	46.2	NA		
Naphthalene	µg/kg	2,100.0	34.4	33.4	42.5		
p,p'-DDE	µg/kg		25.9	38.3	3.8		
Phenanthrene	µg/kg	1,500.0	267.5	275.9	310.6		
Pyrene	µg/kg	2,600.0	534.8	1,000	480.0		
Total Chlordane	µg/kg	NA	17.2	23.1	4.0	6.0	
Total DDTs	µg/kg	46.1	49.3	60	13.1		25.4
Total PAHs	µg/kg					1,800.0*	1,800.0*
Total PCBs	µg/kg	180.0	111.5	125.4	21.3	400.0	0.47
Tributyltin	µg/kg		202.0	308	30.0		

* µg/g organic carbon basis

Draft Report for Review Purposes Only

Table 2. Logistic Regression parameters for the regional and national models compared in this study. National values were taken from Field *et al.*, 2002. B0=intercept; B1=slope; T50=calculated concentration corresponding to a toxicity probability of 0.5. Concentrations are on a dry weight basis.

Chemical	Units	LRM			CA LRM			SoCA LRM			NorCA LRM		
		B0	B1	T50	B0	B1	T50	B0	B1	T50	B0	B1	T50
Cadmium	mg/kg	-0.34	2.51	1.4	0.29	3.18	0.8	0.29	3.18	0.81	1.54	3.43	0.36
Copper	mg/kg				-5.59	2.59	145	-6.76	2.78	268	-6.58	3.84	51
Lead	mg/kg	-5.45	2.77	94	-4.72	2.84	46	-8.64	4.82	62			
Mercury	mg/kg				-0.06	2.68	1.1				1.65	3.05	0.29
Nickel	mg/kg							-8.46	5.70	30			
Zinc	mg/kg	-7.98	3.34	245	-5.13	2.42	132	-9.95	4.20	234	-13.77	6.88	100
1-Methylnaphthalene	µg/kg	-4.14	2.10	94									
1-Methylphenanthrene	µg/kg	-3.59	1.75	112									
2,6-Dimethylnaphthalene	µg/kg	-4.05	1.90	133									
2-Methylnaphthalene	µg/kg	-3.76	1.78	128									
Acenaphthene	µg/kg	-3.62	1.75	116									
Acenaphthylene	µg/kg	-2.96	1.38	140									
Benzo(a)pyrene	µg/kg										-2.27	1.19	80
Benzo(b)fluoranthene	µg/kg	-4.54	1.49	1107							-4.56	2.33	90
Biphenyl	µg/kg	-4.11	2.21	73									
Chlordane, alpha-	µg/kg				-3.41	4.46	5.8	-3.41	4.46	5.8			
Chlordane, gamma-	µg/kg							-3.64	4.18	7.4			
Chrysene	µg/kg										-2.54	1.28	95
Dieldrin	µg/kg	-1.17	2.56	2.9	-1.83	2.59	5.1	-1.24	4.25	2.0			
Fluoranthene	µg/kg	-4.46	1.48	1034									
Fluorene	µg/kg	-3.71	1.81	114									
HMW PAH	µg/kg				-8.19	2.00	12506	-8.19	2.00	12506	-4.26	1.47	785.2
LMW PAH	µg/kg				-6.81	1.88	4127	-6.81	1.88	4127	-3.37	1.49	185.2
Naphthalene	µg/kg	-3.78	1.62	217									
Nonachlor trans	µg/kg				-4.26	5.31	6.3	-4.26	5.31	6.3			
o,p'-DDD	µg/kg							-2.01	3.29	4.1	1.07	2.01	0.3
p,p'-DDD	µg/kg	-1.90	1.49	19				-1.76	2.00	7.6	-0.76	2.45	2.0
p,p'-DDT	µg/kg				-3.55	3.26	12	-1.45	1.60	8.1	-0.55	3.31	1.5
Phenanthrene	µg/kg	-4.46	1.68	455									
Total DDTs	µg/kg										-1.33	2.75	3.0
Total PCBs	µg/kg	-3.46	1.35	368	-4.41	1.48	945	-4.41	1.48	945	-4.41	1.48	945

Draft Report for Review Purposes Only

Table 3. Distribution of sediment chemistry data for the California samples used in the analysis.

Chemical	Units	Northern California			Southern California		
		N	50 th	90 th	N	50 th	90 th
			Percentile	Percentile		Percentile	Percentile
2-Methylnaphthalene	µg/kg	367	10.6	27.2	713	9.6	49.1
Acenaphthene	µg/kg	407	6.0	21.2	674	5.1	46.0
Acenaphthylene	µg/kg	398	8.2	24.3	671	6.2	79.0
Anthracene	µg/kg	422	20.2	91.1	771	18.0	370
Arsenic	mg/kg	393	8.5	12.9	828	8.6	17.3
Benz(a)anthracene	µg/kg	427	63.8	189	838	44.9	720
Benzo(a)pyrene	µg/kg	430	95.7	289	845	65.9	1100
Cadmium	mg/kg	420	0.2	0.4	850	0.4	1.4
Chlordanes, total	µg/kg	404	0.8	3.3	816	7.1	34.3
Chromium	mg/kg	329	122	245	851	56	95
Chrysene	µg/kg	427	72	229	847	64	1090
Copper	mg/kg	405	40.1	65.5	851	76.5	252
DDTs, total	µg/kg	404	3.6	12.4	816	21.4	112
Dibenz(a,h)anthracene	µg/kg	412	12.1	32.5	787	19.1	230
Dieldrin	µg/kg	368	0.2	0.9	297	1.0	3.4
Fluoranthene	µg/kg	425	151	423	849	89.9	1320
Fluorene	µg/kg	414	9.3	34.4	708	6.9	77.5
Lead	mg/kg	409	21.2	37.8	851	35.9	101
Mercury	mg/kg	430	0.3	0.4	843	0.2	0.9
Naphthalene	µg/kg	365	20.9	51.2	733	9.4	44.3
Nickel	mg/kg	399	84.0	114.6	838	20.7	36.6
PCB, total	µg/kg	351	7.9	32.0	851	24.8	196.2
Phenanthrene	µg/kg	392	75.4	242	815	39.8	429
Pyrene	µg/kg	427	190	520	850	102	1500
Silver	mg/kg	418	0.2	0.5	839	0.4	1.4
PAHs, total	µg/kg	431	945	2492	851	619	8573
Zinc	mg/kg	409	110	164	851	180	369

Draft Report for Review Purposes Only

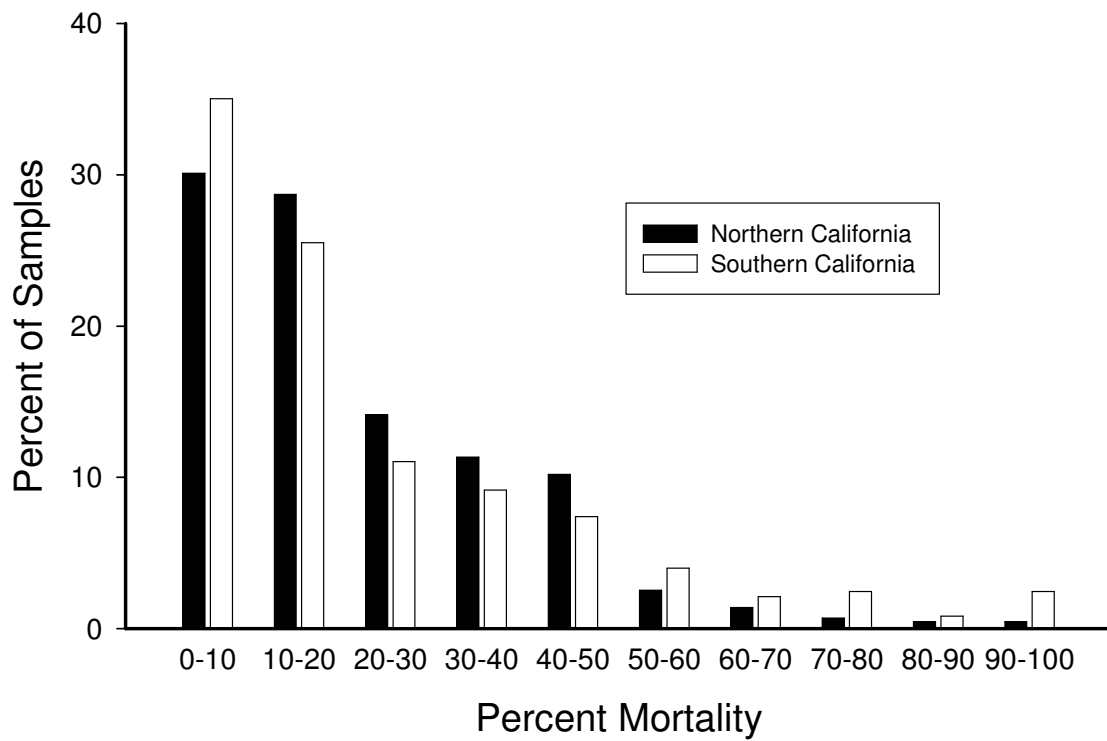


Figure 1. Distribution of sediment toxicity data (10-day amphipod mortality) for the California samples used in the analysis.

Draft Report for Review Purposes Only

Table 4. Thresholds used for evaluations of SGQ index classification accuracy. Nontoxic: <Low threshold; Low Toxicity: Low threshold - <Moderate threshold; Moderate Toxicity: Moderate threshold - <High threshold; High toxicity: >High threshold.

SQG Approach	Index	Low Threshold			Moderate Threshold			High Threshold		
		North	South	State	North	South	State	North	South	State
National ERM	Mean Quotient	0.08	0.06	0.07	0.15	0.12	0.13	0.29	0.38	0.33
National LRM	Maximum Probability	0.17	0.23	0.20	0.26	0.44	0.35	0.50	0.61	0.55
Consensus	Mean Quotient	0.15	0.14	0.14	0.23	0.26	0.25	0.51	0.60	0.55
SQGQ1	Mean Quotient	0.06	0.16	0.10	0.11	0.34	0.19	0.33	0.80	0.52
CA LRM	Maximum Probability	0.25	0.42	0.34	0.42	0.58	0.50	0.62	0.72	0.67
CA ERM	Mean Quotient	0.15	0.14	0.15	0.23	0.25	0.24	0.68	1.28	0.93

Draft Report for Review Purposes Only

Table 5. Nonparametric Spearman correlation (r) and classification accuracy of statewide SQG approaches with amphipod mortality. Values in the shaded cells are within the 90th percentile of the highest median value for the bootstrapped analyses. Analyses were conducted on the combined data for the north and south validation data sets and used thresholds developed using the statewide data set.

Region	Approach	Weighted Kappa	% Agreement	r
State	CA LRM	0.23	37	0.35
State	National ERM	0.17	32	0.25
State	Consensus	0.17	31	0.25
State	National LRM	0.15	35	0.22
State	CA ERM	0.17	33	0.20
State	SQGQ1	0.12	32	0.16

Draft Report for Review Purposes Only

Table 6. Classification accuracy and Spearman correlation of SQG approaches applied to data from each region separately. Values in the shaded cells are within the 90th percentile of the highest median value of the bootstrapped analyses. Analyses were conducted separately using thresholds developed with statewide and region-specific data sets.

Approach	Northern California			Southern California		
	Weighted Kappa	% Agreement	r	Weighted Kappa	% Agreement	r
Statewide Thresholds						
CA LRM	0.20	38	0.39	0.25	35	0.42
National ERM	0.12	27	0.31	0.21	38	0.28
Consensus	0.12	28	0.23	0.22	36	0.31
National LRM	0.11	35	0.18	0.18	34	0.33
CA ERM	0.21	33	0.22	0.15	34	0.18
SQGQ1	0.13	35	0.25	0.10	28	0.26
Region-specific Thresholds						
CA LRM	0.16	27	0.39	0.28	40	0.42
National ERM	0.17	30	0.31	0.22	38	0.28
Consensus	0.15	29	0.23	0.25	39	0.31
National LRM	0.20	33	0.15	0.22	36	0.33
CA ERM	0.21	33	0.22	0.13	33	0.18
SQGQ1	0.21	33	0.25	0.18	33	0.26
Nor/SoCA LRM	0.21	33	0.27	0.22	36	0.37
Nor/SoCA ERM	0.20	35	0.22	0.18	35	0.18